

# THE BEHAVIOR OF BALANCED HALF-SAMPLE VARIANCE ESTIMATES FOR LINEAR AND COMBINED RATIO ESTIMATES WHEN STRATA ARE PAIRED TO FORM PSEUDO-STRATA

Edward J. Stanek III and Stanley Lemeshow, University of Massachusetts, Amherst

## Summary

Expressions are developed illustrating the effect that pairing of strata into pseudo-strata has on balanced half-sample estimates of the variance for estimates of the stratified mean. The development is extended to variance estimates of the combined ratio estimate by using a Monte Carlo sampling experiment. An evaluation is then made of the effect that pairing strata into pseudo-strata has on variance estimates for heights and weights from Cycle 2 of the Health Examination Survey.

The results of this investigation demonstrate that in certain situations, pairing of strata into pseudo-strata can result in highly variable and biased variance estimates of linear and combined ratio estimates. Variance estimates of heights and weights in Cycle 2 of the Health Examination Surveys were found to be insensitive to different pairings of strata, regardless of whether or not the pairings were done in a homogeneous fashion.

## 1. Introduction

Variance estimation has been a problem in the past when dealing with large, complex surveys. Exact expressions for the variance of parameter estimates in such surveys are often unknown and intractable. The balanced half-sample method has grown to be a popular method of estimating variances and is currently used by the National Center for Health Statistics (NCHS) in its Health Examination Survey (HES) and Health Interview Survey (HIS).

As originally developed by McCarthy (1966), the balanced half-sample method requires what may be thought of as the simplest of all designs: a simple stratified sample with two independent observations per stratum. Often, large scale sample surveys are designed so that one cluster of observations (PSU) is selected from each stratum. In these surveys, in order to conform to the balanced half-sample method, strata were paired forming the "pseudo-strata" that were used in the subsequent variance estimation. There were two PSU's per pseudo-stratum. These two PSU's were treated as the two independent observations required by the balanced half-sample variance estimation method.

Properties of the balanced half-sample method have been investigated by numerous authors for a variety of estimates (McCarthy (1966, 1969), Kish and Frankel (1968, 1970), Frankel (1971), Bean (1975), Lemeshow and Epp (1977), Lemeshow and Levy (1977), Lemeshow and Stanek (1977)). All of these studies have started with the assumption upon which the balanced half-sample technique is based: there are two independent observations per stratum. Kish and Frankel (1968, p. 21) suggest that "a model of two

independent primary selections per stratum is probably the most basic design that conforms adequately." Nevertheless, all of the investigators are cognizant of the fact that in situations where the balanced half-sample technique is being applied, the assumption of two independent observations per stratum is seldom met. There is no documented evidence that the techniques used in forming pseudo-strata in the HES produce observations which conform adequately to the assumption of independent observations.

This paper investigates the effect pairing of strata into pseudo-strata has on estimates of the variance for a linear stratified mean, and a combined ratio estimate. The results of this investigation are then placed in the context of height and weight measurements made on children in Cycle 2 of the HES.

## 2. The Linear Case

Consider a situation with  $2L$  strata having means  $\mu_j$ ,  $j=1, \dots, 2L$  and common variance  $\sigma^2$ . As shown by Stanek (1977), for a particular arrangement of the  $2L$  strata into  $L$  pseudo-strata, the expected value of the variance estimate is given by

$$E[\hat{\text{var}}(\bar{x}_{st})] = \frac{\sigma^2}{2L} + \frac{1}{4L^2} \sum_{i=1}^L (\mu_{ij} - \mu_{ij'})^2 \quad (2.1)$$

and the variance of this estimate is given by

$$\text{var}[\hat{\text{var}}(\bar{x}_{st})] = \frac{\sigma^4}{2L^3} + \frac{\sigma^2}{2L^4} \sum_{i=1}^L (\mu_{ij} - \mu_{ij'})^2 \quad (2.2)$$

where  $j$  and  $j'$  represent the two original strata which were combined to form the  $i$ th pseudo-stratum. In these expressions, we appeal to the fact that for linear estimates, the balanced half-sample method produces variance estimates identical to usual stratified sampling formulae. Both Kish (1965, p. 283) and Cochran (1963, p. 141) have noted similar effects in a slightly different context.

Clearly, this process of forming pseudo-strata has certain inherent dangers. Normally, if many of the pseudo-strata in a particular arrangement were comprised of heterogeneous strata, the resulting estimates of variance could be highly biased and variable.

## 3. The Combined Ratio Estimate

Through the use of a sampling experiment, Stanek (1977) investigated the effect that pairing strata to form pseudo-strata has on estimates of variance for combined ratio estimates. The balanced half-sample estimate considered was defined to be

$$\hat{V}(r) = \frac{1}{M} \sum_{i=1}^M (\hat{r}^{(m)} - \bar{r})^2 \quad (3.1)$$

where  $\hat{r}^{(m)}$  is the estimate of the combined ratio estimate for the  $m$ th half-sample and  $\bar{r}$  is the estimate obtained using all the sample observations.

Specifically, let us now assume we have 6 strata with two pairs of observations per stratum.

Strata					
1	2	3	...	6	
$(x_{11}, y_{11})$	$(x_{21}, y_{21})$	$(x_{31}, y_{31})$	...	$(x_{61}, y_{61})$	
$(x_{12}, y_{12})$	$(x_{22}, y_{22})$	$(x_{32}, y_{32})$	...	$(x_{62}, y_{62})$	

The combined ratio estimate of  $R = Y/X$  is

$$\hat{R} = \frac{\sum_{i=1}^6 w_i \bar{y}_i}{\sum_{i=1}^6 w_i \bar{x}_i}$$

$$\text{where } \bar{y}_i = \frac{2}{\sum_{j=1}^2 y_{ij}} \text{ and } \bar{x}_i = \frac{2}{\sum_{j=1}^2 x_{ij}}.$$

Throughout the sampling experiment, we assume that we have normally distributed strata of equal weights. We restrict ourselves to a situation where for each stratum,  $Y = R^{(i)}X$ , i.e., the regression line of  $Y$  on  $X$  for each stratum passes through the origin. With this restriction  $\hat{R}$  becomes an unbiased estimate of the ratio  $R$ . We will assume the correlation coefficient across all strata is constant and equal to  $\rho=0.9$ . Without loss of generality, the mean of  $X$  is held constant across the strata. We will also assume that the variance of  $X$  is held constant across the strata. The four values of the variance of  $X$  as considered in the sampling experiment were  $\sigma_X^2=0.1, 0.001, 0.001, 0.0001$ .

The strata were generated so that the first and second strata had the same set of population parameters with ratio  $R^{(1)}$ , the third and fourth strata had the same set of population parameters different from the first with ratio  $R^{(2)}=dR^{(1)}$ , the fifth and sixth strata has the same set of population parameters different from the previous four with ratio  $R^{(3)}=dR^{(2)}$ . A balanced half-sample estimate of the variance of the combined ratio estimate was made on the basis of all 6 strata. The 6 strata were then paired in all possible ways (i.e.,  $\frac{6!}{2^{3!}} = 15$ ) to form 3 pseudo-strata. Since pairs of strata had the same parameters, these 15 different pairings represented only 5 distinctly different situations. If we assume the strata were paired at random to form pseudo-strata, each of the 15 arrangements would be equally likely. The 5 distinct arrangements of the pseudo-strata along with their probability of occurrence are given below.

$k^{\text{th}}$ Arrangement	Prob (k)	Pseudo- strata 1		Pseudo- strata 2		Pseudo- strata 3	
1	1/15	A	A	B	B	C	C
2	2/15	A	B	A	B	C	C
3	2/15	A	C	B	B	A	C
4	2/15	A	A	B	C	B	C
5	8/15	A	B	A	C	B	C

To obtain a variance estimate using one arrangement of the pseudo-strata, each pair of observations generated per stratum was averaged and then considered as an observation for the pseudo-strata. In such a manner, two pairs of observations were created for each pseudo-stratum from the original data. A balanced half-sample estimate of the variance of the combined ratio estimate was then made on the basis of the 3 pseudo-strata. The constant "d" was chosen such that  $\sum_{i=1}^3 R^{(i)}=5$ . This restriction kept the true ratio  $R$  constant throughout the sampling experiment. The values of  $d$  considered along with the corresponding values of  $R^{(1)}$  are given below.

$d$	$R^{(1)}$
1	5
1.001	4.9950
1.050	4.7581
1.100	4.5317
1.200	4.1209

Balanced half-sample variance estimates produced from the 5 distinct arrangements of pseudo-strata were compared to estimates obtained based on all six strata.

The results of the sampling experiment closely followed results that were derived in the linear case. As the variance of  $X$  became smaller, the relative-bias of variance estimates became larger for a given pairing of the strata. As strata were paired more heterogeneously, the rel-bias of balanced half-sample variance estimates increased. As the difference in strata ratios became larger, balanced half-sample variance estimates became more biased. These results are presented in Table 1.

Similar results held for the change in the variance of the variance estimates (Stanek, 1977).

#### 4. Application to the HES

The increased bias and variability of balanced half-sample variance estimates based on pseudo-strata is only of interest to the extent that surveys actually use this variance estimation method in practice. As has been mentioned earlier, the Health Examination Survey has paired strata into pseudo-strata and used the balanced half-sample method of variance estimation in the past. The Health Interview Survey is presently pairing strata into pseudo-strata and using the balanced half-sample method to estimate the variance. This section will consider

the effect that pairing of strata has on variance estimates for data collected in Cycle 2 of the HES. Specifically, we will investigate the effect that pairing of strata has on variance estimates of height and weight for white children from the ages of 6 to 11. We will use as a reference, variance estimates published on height and weight from the HES. (NCHS, 1972, p. 42.)

The HES was designed with 40 strata formed "in a manner which maximized the degree of homogeneity within superstrata with respect to population size, geographic proximity, degree of industrialization, and degree of urbanization." (NCHS, 1973, p. 6.) One ultimate cluster of observations was chosen from each stratum, and approximately 180 subsequent observations were taken within the ultimate cluster. An estimate of the effect that pairing of strata has on variance estimates cannot be made through a comparison of balanced half-sample variance estimates based on 40 strata with balanced half-sample variance estimates based on 20 pseudo-strata. Differences in these two estimates may stem from the effects of pairing or from the effects of the covariance of observations within the ultimate cluster. Estimates of the effect of pairing strata into pseudo-strata can be made, however, through a comparison of variance estimates under a number of plausible rearrangements of the strata. It is in this manner that we will assess the effect that the formation of pseudo-strata had on variance estimates for heights and weights of children.

The investigators in the HES were cognizant of the potential dangers in forming pseudo-strata to estimate the variance. An effort was made to pair strata as homogeneously as possible. They were paired "on the basis of (1) some subjective determination of the homogeneity of the population in which the primary considerations were population density, region, rate of growth, and industry and (2) concern that strata of approximately equal size would be paired." (NCHS, 1973, p. 27.)

Population density along with rate of growth was defined on a sliding scale for each of 4 geographic regions. The resulting pairings of strata into pseudo-strata for Cycle 2 of the HES are given in Table 2. Clearly, the specific pairing of strata into pseudo-strata as in Cycle 2 of the HES was not the only possible way of forming homogeneous pseudo-strata. A comparison of variance estimates resulting from the HES's pairing with estimates based on other possible homogeneous arrangements (arrangements A-C) and an extremely heterogeneous arrangement (arrangement D) will give a measure of the effect pairing strata has on variance estimates. Details of the criteria for alternative pairing of strata are given by Stanek (1977). It should be noted that the first pairing given in Table 2 is the same as was used by the HES except that when forming their estimates, the HES divided the self-representing strata (pairs 17 through 20) by segments to form new strata. These new strata were used by the HES as pairs 1 through 4 for variance estimation.

Table 3 presents estimates of the standard error of heights and weights of white boys and girls (6-11 years of age) based on the different rearrangements of strata into pseudo-strata. It is important to note that in this regard, there is no asymptotic variance estimate or target value with which to compare variance estimates based on various arrangements of pseudo-strata. Differences in variance estimates for different arrangements may be due to the heterogeneity of strata composing the pseudo-strata, or due to the random variability of samples selected. If consistently large differences were to occur in variance estimates for different arrangements of pseudo-strata, we would suspect that variance estimates were sensitive to pseudo-strata formation. A comparison of estimates based on these alternative homogeneous pairings of strata into pseudo-strata indicates whether variance estimates are highly sensitive to strata pairing. A comparison of variance estimates made when strata are paired heterogeneously with variance estimates made with the HES's pairing should detect gross effects due to the formation of pseudo-strata.

Table 3, which presents a comparison of the standard error estimates for 12 age-sex categories, demonstrates the insensitivity of estimates to alternative pairings of the strata. In most cases, estimates of the standard error based on different arrangements of pseudo-strata differed from published estimates by less than 25%. Due to the small value of the coefficient of variation, this difference would seldom be of practical significance. The largest differences from published estimates of the standard error occurred for arrangements C and D for height of 10 year old white girls. In those cases, estimates of the standard error differed from published estimates by 112%.

Estimates of the standard error based on arrangement D of the pseudo-strata were anticipated to be larger than estimates based on other arrangements. The hypothesis that standard error estimates based on arrangement D of the pseudo-strata were equal to standard error estimates based on another arrangement of pseudo-strata was tested against the one sided alternative that standard error estimates based on arrangement D were greater than standard error estimates based on the other arrangement. The tests were based on Freedman rank sums. (Hollander and Wolfe, 1973, p. 155.) The tests were made on standard error estimates for height and for weight. In each test, the published variance estimates were included as a comparison group. In none of the tests was the hypothesis of equality of standard error estimates rejected in favor of the one sided alternative at  $\alpha=.05$ .

## 5. Conclusions

In summary, balanced half-sample estimates of the variance of mean heights and weights of children were not found to be highly dependent on the arrangement of pseudo-strata for the specific age-color-sex classes considered from Cycle 2 of the HES. Differences did occur due to the arrangement of strata into pseudo-strata but

these differences were no greater than were found by using the complements of the appropriate Plackett-Burman matrices. (See Stanek, (1977).) Rarely did an alternative estimate of the standard error exceed twice the published estimate. Since the coefficient of variation for heights and weights was extremely small for these measurements, differences in estimates of the standard error may not be of practical significance. Only a limited number of situations were considered using HES data. Sampling experiment results demonstrated that in certain situations, large biases could be introduced through the formation of pseudo-strata. Variables whose sample measurements differ widely from stratum to stratum will be more susceptible to these biases. Caution should be exercised in using such variance estimates. Care should be taken to avoid the design of surveys which face this problem in the future.

#### Acknowledgements

This research is based, in part, on the thesis of the first author at the University of Massachusetts, Amherst. Computing assistance was obtained from the University Computing Center at the University of Massachusetts.

#### Bibliography

- Bean, Judy A. (1975). "Distribution and Properties of Variance Estimators for Complex Multistage Probability Samples: An Empirical Distribution." Data Evaluation and Methods Research. NCHS, Series 2, #65. DHEW Publication No. (HRA) 75-1339.
- Cochran, W.G. (1962). Sampling Techniques. New York: John Wiley and Sons.
- Hollander, M. and Wolfe, D.A. (1973). Non Parametric Statistical Methods. New York: John Wiley and Sons.
- Kish, Leslie (1965). Survey Sampling. New York: John Wiley and Sons.
- Kish, Leslie and Frankel, Martin R. (1968). "Balanced Repeated Replication for Analytical Statistics." Proceedings of the Social Statistics Section of American Statistical Association: 2-10.
- Kish, Leslie and Frankel, Martin R. (1970). "Balanced Repeated Replication for Standard Errors." Journal of the American Statistical Association, Vol. 65, No. 331, pp. 1071-1094.
- Lemeshow, S. and Epp, R. (1977). "Properties of the Balanced Half-Sample and Jackknife Variance Estimation Techniques in the Linear Case." Communications in Statistics A, Vol. 6, Issue 13.
- Lemeshow, S. and Levy, P. (1977). "Estimating the Variance of Ratio Estimates in Complex Sample Surveys with Two Primary Sampling Units Per Stratum - A Comparison of Balanced Replication and Jackknife Techniques." Submitted for publication.
- Lemeshow, Stanley and Stanek, Edward J. (1977). "Estimating the Variance of the Slope of a Linear Regression in a Stratified Random Sample with the Balanced Half-Sample Technique." Submitted for publication.
- McCarthy, Philip J. (1966). "Replication. An Approach to the Analysis of Data from Complex Surveys." Vital and Health Statistics. NCHS, Series 2, #14.
- McCarthy, Philip J. (1969). "Pseudoreplication: Further Evaluation and Application of the Balanced Half-Sample Technique." Vital and Health Statistics. NCHS, Series 2, #31.
- National Center for Health Statistics (1972). "Height and Weight of Children: Socio-economics Status, United States." Vital and Health Statistics. NCHS, Series 11, #119.
- National Center for Health Statistics (1973). "Sample Design and Estimation Procedures for a National Health Examination Survey of Children." Vital and Health-Statistics. NCHS, Series 2, #43.
- Plackett, R.L. and Burman, J.P. (1946). "The Design of Optimum Multifactorial Experiments." Biometrika, Vol. 33 (pt. IV): pp. 305-325.
- Stanek, E.J. (1977). "The Properties of Balanced Half-Sample Variance Estimates in Complex Surveys when Strata are Paired to Form Pseudo-Strata." Biostatistics Technical Reports, No. 77-8, Division of Public Health, School of Health Sciences, University of Massachusetts/Amherst.

TABLE 1

Results of a sampling experiment for  $\hat{V}_2(\hat{R})$  in which  $N=2$  observations were selected from each of  $L=6$  strata. The strata were paired in  $k=1, \dots, 5$  arrangements of pseudo-strata. Six multiplicative factors,  $d$ , (1, 1.001, 1.010, 1.050, 1.100, and 1.200) and four rel-variances of  $X$ , (0.1, 0.01, 0.001, 0.0001) were used. In each case,  $\rho=0.9$ .

		Rel-Bias						
Rel-Var (X)	k	1	1.001	1.010	d	1.050	1.100	1.200
.1000	1	-.926E-01	-.924E-01	-.909E-01	-.842E-01	-.760E-01	-.609E-01	
.1000	2	-.646E-01	-.642E-01	-.583E-01	.110E-01	.173E+00	.618E+00	
.1000	3	-.537E-01	-.528E-01	-.361E-01	.225E+00	.918E+00	.310E+01	
.1000	4	-.361E-01	-.361E-01	-.346E-01	.217E-01	.199E+00	.843E+00	
.1000	5	-.513E-01	-.510E-01	-.416E-01	.139E+00	.639E+00	.224E+01	
.0100	1	-.586E-01	-.585E-01	-.568E-01	-.501E-01	-.427E-01	-.318E-01	
.0100	2	-.541E-01	-.525E-01	-.143E-01	.618E+00	.222E+01	.675E+01	
.0100	3	-.446E-01	-.419E-01	.791E-01	.261E+01	.969E+01	.324E+02	
.0100	4	-.405E-02	-.440E-02	.171E-01	.642E+00	.257E+01	.947E+01	
.0100	5	-.257E-01	-.241E-01	.635E-01	.194E+01	.723E+01	.243E+02	
.0010	1	-.502E-01	-.501E-01	-.484E-01	-.415E-01	-.342E-01	-.239E-01	
.0010	2	-.508E-01	-.440E-01	.259E+00	.631E+01	.222E+02	.674E+02	
.0010	3	-.424E-01	-.273E-01	.109E+01	.261E+02	.969E+02	.325E+03	
.0010	4	.730E-02	.793E-02	.261E+00	.672E+01	.263E+02	.961E+02	
.0010	5	-.188E-01	-.787E-02	.825E+00	.196E+02	.727E+02	.244E+03	
.0001	1	-.478E-01	-.476E-01	-.459E-01	-.389E-01	-.316E-01	-.216E-01	
.0001	2	-.499E-01	-.963E-02	.278E+01	.624E+02	.220E+03	.672E+03	
.0001	3	-.418E-01	.802E-01	.110E+02	.260E+03	.968E+03	.325E+04	
.0001	4	.110E-01	.317E-01	.270E+01	.679E+02	.265E+03	.965E+03	
.0001	5	-.168E-01	.744E-01	.825E+01	.195E+03	.726E+03	.244E+04	

TABLE 2

Arrangements of strata into pseudo-strata					
	HES	A	B	C	D
Boston, Mass.	1	1	3	3	1
Neward, N.J.	1	4	10	10	2
Jersey City, N.J.	2	4	10	10	3
Allentown, Pa.	2	3	9	6	4
Poughkeepsie, N.Y.	3	5	14	11	5
Hartford, Conn.	3	3	8	6	6
Columbia, S.C.	4	7	11	12	7
Charleston, S.C.	4	8	11	12	8
Marked Tree, Ark.	5	10	17	15	4
Georgetown, Del.	5	10	17	16	9
Barbourville, Ky.	6	9	15	15	10
West Liberty, Ky.	6	9	15	16	3
Cleveland, Ohio	7	12	7	5	11
Minneapolis, Minn.	7	12	8	7	2
Lapeer, Mich.	8	15	18	17	11
Ashtabula, Ohio	8	14	18	18	12
San Francisco, Calif.	9	17	6	5	12
Denver, Colo.	9	18	7	8	13
Prowers, Colo.	10	20	19	20	1
Maripose, Calif.	10	19	19	19	14
Atlanta, Ga.	11	6	4	9	15
Houston, Tex.	11	17	6	8	10
Des Moines, Iowa	12	13	12	14	13
Wichita, Kans.	12	18	13	13	16
Birmingham, Ala.	13	7	9	9	17
Grand Rapids, Mich.	13	13	12	14	18
Clark, Wis.	14	15	20	17	6
Grant, Wash.	14	20	16	19	15
Portland, Maine	15	5	14	11	19
Ottumwa, Iowa	15	14	20	18	17
Sarasota, Fla.	16	8	13	13	20
Brownsville, Tex.	16	19	16	20	20
Philadelphia, Pa.	17	1	3	3	16
Baltimore, Md.	17	6	4	7	18
Chicago, Ill.	18	11	5	4	5
Detroit, Mich.	18	11	5	4	19
Los Angeles, Calif.	19	16	2	2	7
Los Angeles, Calif.	19	16	2	2	14
New York, N.Y.	20	2	1	1	8
New York, N.Y.	20	2	1	1	9

Table 3

Standard error estimates for heights in CM and weights in KG for white boys and girls in 6 age categories from Cycle 2 of the health examination survey. Standard error estimates are presented for each of 5 alternative arrangements of strata into psuedo-strata, (HES, A, B, C, and D). The published st. error-estimates and mean heights are given for comparison. (NCHS, 1972)

	HEIGHT											
	Boys Age						Girls Age					
	6	7	8	9	10	11	6	7	8	9	10	11
Published Mean Height	118	124	130	135	140	146	118	123	129	135	141	147
Published St. Error	.30	.38	.29	.50	.37	.30	.32	.17	.39	.36	.34	.37
HES	.37	.35	.26	.46	.37	.38	.28	.21	.51	.44	.47	.37
A	.37	.35	.26	.44	.23	.35	.31	.26	.33	.45	.67	.45
B	.41	.32	.25	.46	.40	.39	.41	.21	.25	.45	.39	.24
C	.28	.37	.25	.39	.36	.31	.22	.25	.30	.42	.72	.36
D	.29	.36	.31	.46	.36	.33	.31	.22	.32	.37	.72	.29
	WEIGHT											
	Boys Age						Girls Age					
	6	7	8	9	10	11	6	7	8	9	10	11
Published Mean Height	22	25	28	31	34	39	22	24	28	31	35	40
Published St. Error	.17	.21	.25	.47	.30	.40	.25	.20	.26	.43	.44	.36
HES	.19	.22	.28	.46	.31	.46	.25	.18	.28	.48	.47	.42
A	.22	.20	.24	.38	.22	.54	.21	.22	.27	.40	.52	.41
B	.20	.22	.23	.31	.28	.43	.27	.24	.21	.51	.40	.53
C	.15	.19	.28	.37	.29	.36	.25	.26	.24	.45	.64	.44
D	.16	.23	.29	.35	.36	.46	.24	.20	.27	.46	.47	.31